

Fully Convolutional Neural Networks for Raw Eye Tracking Data Segmentation, Generation, and Reconstruction

Wolfgang Fuhl

Human-Computer Interaction
Eberhard Karls University Tübingen
Germany, Tübingen, Sand 13
wolfgang.fuhl@uni-tuebingen.de

Yao Rong

Human-Computer Interaction
Eberhard Karls University Tübingen
Germany, Tübingen, Sand 13
yao.rong@uni-tuebingen.de

Enkelejda Kasneci

Human-Computer Interaction
Eberhard Karls University Tübingen
Germany, Tübingen, Sand 13
enkelejda.kasneci@uni-tuebingen.de

Abstract—In this paper, we use fully convolutional neural networks for the semantic segmentation of eye tracking data. We also use these networks for reconstruction, and in conjunction with a variational auto-encoder to generate eye movement data. The first improvement of our approach is that no input window is necessary, due to the use of fully convolutional networks and therefore any input size can be processed directly. The second improvement is that the used and generated data is raw eye tracking data (position X, Y and time) without preprocessing. This is achieved by pre-initializing the filters in the first layer and by building the input tensor along the z axis. We evaluated our approach on three publicly available datasets and compare the results to the state of the art.

I. INTRODUCTION

Eye movements are the basis to get more information about a person. Most research papers investigate intentions, cognitive states [1], workload [2] and attention of a person. The eye movements are used to generate more complex features for machine learning to classify or regress the desired information. This knowledge about a person is important in multiple fields, like automated driving [3] and for measuring the work load of a surgeon [4]. For the eye movements themselves, there are also application areas like the recognition of eye diseases [5] and the foveated rendering [6]. The fields of eye tracking applications are becoming more and more diverse, but even today there are still a multitude of challenges. One of these challenges is the reliable classification of eye movements based on raw data. The commonly used algorithms require the determination of a large number of thresholds [7]. But most algorithms are bound to certain sampling rates of the eye tracker and do not work even if the signal is very noisy [8]. Newer approaches avoid these limitations by using machine learning methods. This allows the algorithm to be re-trained for any eye tracker. The preprocessing of the data is still used by these methods; however, it brings restrictions regarding data in which the preprocessing does not work as intended. Another problem of machine learning is the necessity of annotated data. For this purpose simulators have already been presented [9], [10] that address this challenge.

In this paper, we present an approach that is not bound to a window size. We achieve this by the exclusive use of convolution layers that are spatial invariant and not bounded to an input size. Compared to other machine learning approaches, our approach uses raw data as input, eliminating preprocessing. This has the advantage that our approach works autonomously and does not depend on the effectiveness of other methods. Furthermore, we show that our approach can be used for the classification, generation and reconstruction of eye tracking data.

Contribution of this work:

- 1 Processing of raw eye tracking data with neural networks by a sign-based weight pre-initialisation and data arrangement.
- 2 Window free approach (fully convolutional).
- 3 Use of neural networks for eye tracking data reconstruction.
- 4 Use of variational autoencoders to generate eye tracking data.

II. RELATED WORK

Since this work proposes an approach for the classification of eye movements as well as the reconstruction and generation, we have divided this section into two subsections.

A. Eye Movement Classification

The two most famous and most common algorithms in the field of eye movement classification are Identification by Dispersion Threshold (IDT) [11] and Identification by Velocity Threshold (IVT) [11]. In the former, the data is first reduced [12]. Then, two thresholds are used to distinguish between fixations and saccades. The first threshold limits the dispersion of the measurement points and the second threshold limits the minimum duration of fixations. In the second algorithm, however, only one threshold is used, which limits the eye movement velocity. If an eye movement is above this threshold, it is classified as Saccade, otherwise, a fixation is assumed. This second algorithm (IVT) has already been extended by adaptive methods to determine the

threshold [13]. A further improvement in the signal noise level adaptation was achieved by using the Kalman filter (IKF) [14]. Here, the Kalman filter is used to predict the next value, resulting in the signal being smoothed online. In addition to the velocity threshold, a threshold is used for the minimum fixation duration. A similar algorithm has been published in [15]. The difference to IKF is the use of the χ^2 -test instead of the Kalman filter. Not only has the IVT algorithm been extended, but also the IDT algorithm. The first extension is the F-tests Dispersion Algorithm (FDT) [16]. The F-test provides the probability whether several data points belong to the same class. Since the F-test always expects a normal distribution, it is relatively susceptible to noise in the data. In the Covariance Dispersion Algorithm (CDT) [17], the F-test was replaced by a covariance matrix. For classification, the CDT requires three thresholds. The first two thresholds are for the variance and the covariance and thus represent an improvement of the dispersion threshold. The third threshold is for the minimum fixation duration. The last approach that followed the idea of IDT is the Identification by a Minimal Spanning Tree (IMST) [15] algorithm. Here a tree structure is calculated on the data, where each data point represents a leaf of the tree. Clusters are formed over the number of branches, which represent fixations and can be seen as a form of dispersion.

The first approach with machine learning was made in for the adaptive thresholds. Hidden Markov Models (HMM) were used to determine the class based on the velocity and the current state of the model [15]. Most models have two states (fixation and saccade) to classify the data. This approach of HMMs was also extended with smooth pursuits [18] and an additional state. In addition to the smooth pursuits, the post saccadic movements (PSM) also became interesting for science. The first algorithm dealing with the detection of PSM was presented in [19]. One year later the Binocular-Individual Threshold (BIT) [20] algorithm was introduced, which uses both eyes to detect PSM. This algorithm also used adaptive thresholds and follows the idea that both eyes perform the same movement. The first algorithm able to detect four eye movements was presented in [21]. This algorithm uses different data cleansing techniques and adaptive thresholds. For eye tracking data with a very high sampling rate an algorithm was presented in [22]. This algorithm is able to detect four eye movement types and is based on several steps in which all data is processed. The first step generates a rough segmentation and, in the following steps, this segmentation is further refined. Meaning, that the algorithm cannot be used online.

Novel approaches for eye movement classification use modern machine learning approaches. The first approach to be mentioned here is [23]. This approach uses a conventional neural net with convolution layers and a fixed window size. The data in each window is first transferred to the frequency domain via the fast Fourier transformation and then used as input for the neural net, which classifies the eye movement type. Another approach is described in [24]. Here, a random

forest is used to be applicable to bending eye tracker sampling rates. For this the input data is interpolated via cubic splines and 14 different features, like the eye movement speed, are calculated. These 14 features serve as input data for the random forest and must always be calculated in advance. In addition, postprocessing is performed with Gaussian smoothing of the class probabilities as well as a heuristic for the final classification. A rule based learning algorithm was presented in [25]. Different data streams like the eye movement speed can be provided to the algorithm, whereupon the algorithm learns rule sets consisting of thresholds. Based on these rule sets, new data is classified to eye movement types. The last representative of the modern methods, is a feature which enters the velocities based on their direction into a histogram [26]. This histogram is normalized and can be used with any machine learning method.

B. Eye Tracking Data Reconstruction & Generation

The synthesis of eye movements is still a challenging task today. The Kalman filter makes the prediction of the next gaze point and is able to generate fixations, saccades and smooth pursuits. A disadvantage of this method is that there is no realistic noise in the signal which reflects the inaccuracy of the eye tracker. A rendering based approach was introduced in 2002 [27]. The main focus of this method was on the saccades but the method is also able to generate smooth pursuits and binocular rotations (vergence). A pure data based approach was presented in [28]. These methods simulate eye movements as well as head rotations. A disadvantage of the methods is that the head movements automatically trigger eye movements. Normally, head movements are only triggered when a target is more than $\approx 30^\circ$ apart [29]. Another purely data-driven approach is described in [30]. It is an automated framework that simulates head, eye and eyelid movements. The method uses sound input to generate the movements, which are projected over several normal distributions onto eye, head and eyelid movements. Another approach focused on eye rotation is described in [9] and is based on the description of eye muscles by [31]. A disadvantage of this simulation is that eye movements cannot be generated automatically, but have to be predefined. The last rendering based approach is described in [32]. Here images and gaze vectors are randomly generated and the simulator is used to train machine learning techniques for detection and gaze vector regression. The methods mentioned so far originate from computer graphics and do not have the aim to generate realistic eye movement sequences. Their actual use lies in the interaction with humans [33]. This results in all movements being error free and absolutely accurate, which does not correspond to reality. Furthermore, the procedures described do not include evaluation of visual input or task specific behavior. The first approach to realistic simulation of eye tracking data for static images is described in [34]. This approach uses a random sequence of numbers in combination with statistical models and saliency maps to generate eye tracking data. An extension of this approach that added noise is described in [35]. In addition to noise, jitter

based on a normal distribution was added [36]. A multi-layer calculation approach is described in [10]. The simulator allows to generate a random sequence of eye movement sequences and to map them to static, dynamic, and eye tracking data. This simulator can also generate any sampling rate as well, as it supports dynamic sampling rates.

Machine learning based approaches have already been presented. In [37] deep recurrent neural networks are used to generate eye movements based on static images. A disadvantage of this approach is that it only works on already seen images. An approach which uses Generative Adversarial Networks (GANs) is described in [38]. This approach uses recurrent layer and a combination of static image and saliency maps to predict a scanpath.

III. METHOD

In this section, we describe our three approaches and how we trained them. Each task (semantic segmentation, reconstruction and generation) has its own subsection and is described in detail together with the training parameters. All networks were trained from scratch with a random initialization. While all of our models work with raw eye tracking data, it has to be mentioned that NaN or Inf values in the input files will corrupt the result. For the reconstruction model, those values have to be set to zero for example.

A. Semantic Segmentation

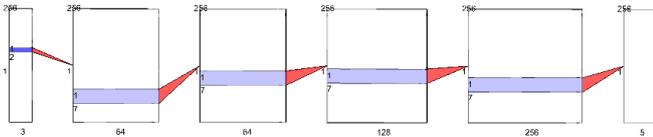


Fig. 1. The eye movement segmentation model used in our experiments.

Our eye movement segmentation network consists of five convolution layers with rectifier linear unit (ReLU) activation units. The input to our network is raw eye tracking data (position x , y and time). For our model, the input data is arranged one after the other (see Figure 1). This results in an input tensor that has a fixed depth of 3, a width of 1, and an arbitrary height. In Figure 1, the height was set to 256. This arrangement has the advantage that the weight tensor of the convolution extends over the whole depth and is only shifted along the height. Therefore, a convolution always sees all three input types (position x , y and time). In our training we used a fixed constant of one hundred as divisor for the input values to gain numerical stability. Without this divisor, it is also possible to train the network, but with lower learning rates, which prolong the whole training.

The first convolution layer has a height of two (see Figure 1). For this layer, it is important to check that, for each superimposed weight (along the height), one is positive and one is negative. Meaning, that after random value initialization, if two superimposed weights in the first layer are both positive, one is set to its negative value and vice versa. This has only

to be done for the first layer. All the other layers are randomly initialized without any modification.

The last layer of our model has five output layers. This is due to the use of the softmax loss function and these five layers hold the output probability distributions for the corresponding eye movement types (Fixation, Saccade, Smooth pursuit, PSM, error) and can be extended. In addition, it can be seen that our network does not use any down or upscaling operation.

1) *Semantic Segmentation training parameters:* For training on both datasets we used an initial learning rate of 10^{-2} together with the stochastic gradient descent (SGD) [39] optimizer. The parameters for the optimizer are $weightdecay = 10^{-4}$ and $momentum = 0.9$. For the loss function we used the weighted log multi class loss together with the softmax function. After each five hundred epochs, the learning rate was reduced by 10^{-1} until it reached 10^{-6} when the training was stopped. For data augmentation, we used random jitter that changes the value of a position to up to 2% around its original value. In addition, we shifted the entire input scanpath by a randomly selected value (the same value for all entries). We also used different input sizes where it has to be noted that in one batch, all length were equal because, otherwise, computational problems arise due to the not aligned data.

B. Reconstruction

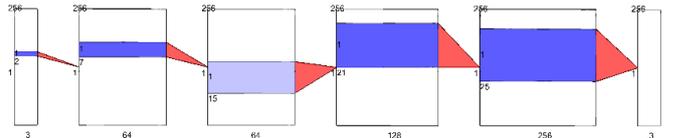


Fig. 2. The model of our eye tracking raw data reconstruction network.

The model we use to reconstruct eye tracking data has the same structure as the segmentation model (Figure 2). The only difference is the output, which corresponds to the eye tracking signal itself. At the beginning, there is the sign-based pre-initialized convolution with the height two. Then follow the convolution layers, where the size of the convolution always doubles after the layer with a convolution height of seven. The last convolution layer reconstructs the signal and has an output depth of three (X , Y , time) and a height of twenty-five.

At this point it must be said that the output as well as the input layer can be extended. An example of this would be three dimensional coordinates, which can be processed and trained with an input and output layer of depth four. Furthermore, as with the segmentation mesh, no input window is required, allowing the mesh to be applied to any input length. Of course it is also possible to train and validate the net with different and varying input lengths.

1) *Reconstruction training parameters:* For training on both datasets, we used an initial learning rate of 10^{-4} and changed it after ten epochs to 10^{-3} . This was done to avoid numerical problems for the random initialized models, which end up in not a number results (NaN). As optimizer, we used adam [40] with the parameters $weightdecay = 5 * 10^{-4}$,

$momentum1 = 0.9$, and $momentum2 = 0.999$. As loss function, we used the L2 loss for the first hundred epochs. Afterwards, we used the L1 loss function to improve the accuracy of the network. The learning rate was decreased by 10^{-1} after each five hundred epochs and the training was stopped at a learning rate of 10^{-6} . For data augmentation, we used random jitter that changes the value of a position to up to 2% from its original value. In addition, we shifted the entire input scanpath by a randomly selected value (the same value for all entries). We also used different input sizes, where it has to be noted that in one batch, all length were equal because otherwise computational problems arise due to the not aligned data. In addition, it is important to note that for the training, only parts without error are selected since otherwise, our network would learn to reconstruct errors or what is most likely, is that it would learn nothing.

C. Generation

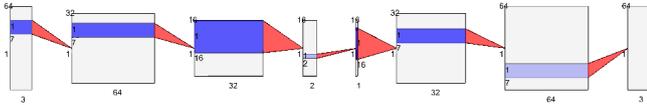


Fig. 3. The Variational Autoencoder (VAE) used for eye tracking data generation.

Figure 3 shows the structure of the variational autoencoder [41] (VAE) used. In comparison to the reconstruction as well as the segmentation net, we did not use the first pre-initialized convolution layer. This is due to the position itself is not used as input nor as output. We used the position change in x and y together with the time difference as input and as target for learning. This was done to make the output dependent to the last position of the scanpath and avoids jumping around the image since a new scanpath is generated based on a number of random values from normal distributions. Therefore, we need to specify randomly an initial start position from where the scanpath is further constructed. The input layer is followed by two convolution layers, which also reduce the input by half. This was realized by average pooling. The last output layer of the encoder has a depth of two and corresponds to the mean value and the variance of the normal distributions. Then, a layer with depth one follows, which corresponds to Z , the value of the learned distribution. The decoder part of the network then learns to generate new data based on the distributions. For completeness, a brief description of the VAE is given below.

1) *Description Variational Autoencoder (VAE)*: A VAE is similar to a normal autoencoder and consists of an encoding and a decoding part. The main difference is that, instead of encoding an input as a single point, the input is encoded as a distribution. Therefore, the encoder learns to map the input to the parameters of the normal distribution (mean m and variance v). The decoder in contrast learns to generate new samples based on the output of the normal distribution z .

Since the error cannot be propagated back through the distribution, the reparametrization trick is used. The calculation

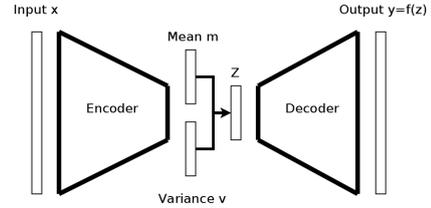


Fig. 4. The concept of an Variational Autoencoder (VAE).

of z ($z = N(m, v)$) is replaced with $z = m + v * N(0, 1)$. This calculation of the distribution is derivable and thus the error can be propagated back. Another difference between the VAE and the normal autoencoder is that the error depends not only on the difference between the input (x) and the output (y) but also on the similarity of the distributions. Therefore, the loss function is added another term, the Kullback-Leibler divergence. This divergence computes the distance between two distributions. The whole loss function is therefore computed as $\|x - y\|^2 - KL(N(m, v), N(0, 1))$.

2) *Generation training parameters*: For training, we used an initial learning rate of 10^{-4} and changed it after one hundred epochs to 10^{-3} . As optimizer, we used stochastic gradient descent (SGD) [39]. The parameters for the optimizer are $weightdecay = 10^{-6}$ and $momentum = 0.9$. As loss function we used the L2 loss in combination with the KL divergence as described in Section III-C1. The learning rate was decreased by 10^{-1} after each thousand epochs and the training was stopped at a learning rate of 10^{-6} . We did not use any data augmentation technique since the reparametrization trick already induces some deformation in the output data.

IV. EVALUATION

The evaluation section is split into three subsections. In each subsection we evaluate our approach for a specific task (semantic segmentation, reconstruction and generation) on multiple publicly available data sets. For evaluation we used the data sets from [18] (DS-SAN) and [8] (DS-AND) where we only used the annotations from MN for the semantic segmentation evaluation. The data set from DS-SAN consists of 24 recordings from six subjects. Each subject made four recordings with different challenges for eye movement detection. The data set contains fixations, saccades, and smooth pursuits. In addition it contains errors from the Dikablis Pro eye tracker, which has a sampling rate of 30 Hz. The subjects were recorded at a distance of 300 mm from the screen using a chin rest. The data set DS-AND contains annotations for fixations, saccades, smooth pursuits, and post saccadic movement (PSM). It was recorded using a SMI HiSpeed 1250 system with a chin and forehead rest. The data set consists of 34 binocular recordings from 17 different students at Lund University and each file has a sampling rate of 500 Hz. They used static and dynamic stimuli during recording.

A. Evaluation Semantic Segmentation

For the comparison of our algorithm to the state of the art we used the algorithms [18] (IBDT), [19] (EV), [42] (I2MC),

TABLE I

DIFFERENT TRAINING CONFIGURATIONS OF THE MACHINE LEARNING APPROACHES USED TOGETHER WITH THE HOV FEATURE [26].

Name	Configuration
knn5-20	$k=5,10,15$, or 20
tree1	Maximum splits 50
tree2	Maximum splits 50, Predictor selection with curvature Exact categorization
tree3	Maximum splits 50, Predictor selection with curvature, Exact categorization, split criterion deviance
tree4	Maximum splits 50, Exact categorization
tree5	Maximum splits 50, Exact categorization, split criterion deviance
svm-lin	Linear kernel function
svm-pol	Second order polynomial as kernel function

[21] (LS), [25] (RULE) , and [26] (HOV). All algorithms where configured for offline use since only three algorithms are configurable for online use (IBDT, HOV, & RULE). In addition, we used the data unmodified, which means that no errors were removed, nor was any preprocessing applied with the exception of preprocessing, which is integrated in the state of the art algorithms. For the evaluation itself, we only considered the annotated data points. For our approach, we used a four fold cross validation where the data of one subject can only be in one fold. The training configurations and machine learning approaches used together with the HOV feature are shown in Table I. As can be seen we applied three well known machine learning approaches namely k nearest neighbors (knn), decision trees (tree), and support vector machines (svm) with different configurations. Table II and Table III show the results for recall ($TP/(TP + FN)$) and precision ($TP/(TP+FP)$). TP are true positives, FP are false positives, and FN are false negatives. Recall therefore stands for the amount of correctly detected eye movement samples. In contrast to this, precision allows to evaluate the reliability of predictions.

Table II shows that our approach outperforms the other state of the art approaches on both data sets for nearly all eye movement types based on pure correct predictions (Recall). It has to be mentioned that we used the raw input data without smoothing out errors as it is done, for example, in IBDT. For Fixations, the HOV feature in combination with the KNN machine learning algorithm performs equal to our fully convolutional network on the DS-SAN data set. However, For the other eye movement types, our approach outperforms the HOV feature with KNN. For the not machine learning based approaches (EV, IBDT, LS, I2MC), it can be seen that they perform well only for the data with a frequency they are designed for. One example for this is IBDT, which performs well on DS-SAN but not on DS-AND due to the higher frequency. In addition, they can of course not detect eye movement types that are not included by the creator of the algorithm. The best example for this is I2MC which can only detect fixations and saccades. Another issue with the data sets itself is that the annotations change especially for saccades. In the DS-SAN data set for example, saccades are annotated

TABLE II

RECALL FOR EACH EYE MOVEMENT TYPE WITH ERRORS IN THE INPUT DATA. PSM STANDS FOR POST SACCADIC MOVEMENT.

Data Alg.	Recall					
	Fixation	Saccade	Pursuit	Noise	PSM	
DS-AND	EV	0.61	0.73	0	0.94	0.02
	IBDT	0.65	0.35	0.63	0	0
	LS	0.91	0.88	0.15	0.13	0
	I2MC	0.02	0.95	0	0	0
	RULE	0.79	0.85	0.69	0.67	0.78
	Proposed	0.94	0.93	0.91	0.96	0.89
DS-SAN	EV	0.18	0.25	0	1.0	-
	IBDT	0.97	0.28	0.84	0	-
	LS	0.95	0	0.06	0	-
	I2MC	0.92	0.10	0	0	-
	RULE	0.94	0.91	0.89	0.65	-
	HOV & knn5	0.97	0.73	0.91	0.70	-
	HOV & knn10	0.98	0.70	0.91	0.70	-
	HOV & knn15	0.98	0.69	0.91	0.70	-
	HOV & knn20	0.98	0.68	0.91	0.70	-
	HOV & tree1	0.97	0.92	0.88	0.72	-
	HOV & tree2	0.97	0.91	0.88	0.72	-
	HOV & tree3	0.97	0.92	0.88	0.73	-
	HOV & tree4	0.97	0.92	0.88	0.72	-
	HOV & tree5	0.97	0.92	0.88	0.77	-
	HOV & svm-lin	0.95	0.85	0.61	0.76	-
	HOV & svm-pol	0.82	0.84	0.90	0.72	-
Proposed	0.98	0.95	0.94	0.89	-	

TABLE III

THE PRECISION FOR EACH EYE MOVEMENT TYPE WITH ERRORS IN THE INPUT DATA. PSM STANDS FOR POST SACCADIC MOVEMENT.

Data Alg.	Precision					
	Fixation	Saccade	Pursuit	Noise	PSM	
DS-AND	EV	0.68	0.37	0	0.73	0.03
	IBDT	0.72	0.70	0.35	0	0
	LS	0.82	0.33	0.63	0.06	0
	I2MC	0.09	0.08	0	0	0
	RULE	0.79	0.65	0.82	0.59	0.32
	Proposed	0.91	0.81	0.90	0.88	0.73
DS-SAN	EV	0.78	0.31	0	0.01	-
	IBDT	0.93	0.73	0.76	0	-
	LS	0.77	0	0.23	0	-
	I2MC	0.76	0.071	0	0	-
	RULE	0.98	0.86	0.89	0.61	-
	HOV & knn5	0.96	0.91	0.91	0.82	-
	HOV & knn10	0.96	0.93	0.92	0.88	-
	HOV & knn15	0.96	0.94	0.91	0.87	-
	HOV & knn20	0.95	0.95	0.91	0.89	-
	HOV & tree1	0.97	0.90	0.89	0.79	-
	HOV & tree2	0.97	0.91	0.89	0.78	-
	HOV & tree3	0.97	0.91	0.89	0.84	-
	HOV & tree4	0.97	0.90	0.89	0.79	-
	HOV & tree5	0.97	0.91	0.88	0.84	-
	HOV & svm-lin	0.91	0.86	0.76	0.79	-
	HOV & svm-pol	0.96	0.42	0.77	0.26	-
Proposed	0.98	0.95	0.94	0.91	-	

after the velocity peek while in DS-AND the velocity profile of the saccade is annotated. This issue makes it impossible to do cross data set evaluations.

Table III shows that our approach outperforms the other state of the art approaches on both data sets for all eye movement types based on the reliability of the predictions (Precision). In combination with Table II, this means that our approach does not only detect a majority of the eye movement types correctly, it is also more reliable in its detections. Since

TABLE IV

THE ABSOLUTE ERROR FOR THE RECONSTRUCTION OF THE X AND Y POSITION AS EUCLIDEAN DISTANCE FOR BOTH DATA SETS. THE ERROR IS UPSCALED TO THE REAL INPUT VALUE RANGE (MULTIPLIED BY 100 TO COMPENSATE FOR THE NORMALIZATION DIVIDER).

Evaluation Data set		Induced Error	Absolut Error
Entire Input Sequence	DS-AND	5%	1.76 px
		10%	2.06 px
		15%	2.37 px
		20%	2.72 px
		25%	3.01 px
	DS-SAN	5%	1.62 px
		10%	1.68 px
		15%	1.69 px
		20%	1.70 px
		25%	1.72 px
Induced Errors Only	DS-AND	5%	19.26 px
		10%	22.37 px
		15%	25.51 px
		20%	28.84 px
		25%	31.68 px
	DS-SAN	5%	7.06 px
		10%	7.25 px
		15%	7.38 px
		20%	7.41 px
		25%	7.46 px
		30%	7.67 px

our approach did not reach 100% for each eye movement type and the variety of challenges in the real world can not easily be covered by scientific data sets, we think that the eye movement detection is still an open problem. The benefits of our approach is the simple realization with modern neuronal network toolboxes. In addition, it is adaptable to new eye movements and varying annotations but this is true for all machine learning based approaches.

B. Evaluation Reconstruction

For the reconstruction, we also used both data sets ([18] & [8]). For the evaluation, a random file was selected one hundred times from the test data set. In this file, a hundred random length and random start positions were selected to extract sections out of the document. In case one section already contained errors, it was discarded and another section was selected. This approach was chosen to evaluate the reconstruction, as our method is not intended to reconstruct an error. To evaluate the reconstruction quality, we injected several fixed percentage amounts of errors for each section. These errors were either the setting of a zero or a random number. Each position for an error injection was selected randomly. As a measure of the quality of the reconstruction, we used the mean absolute error. In addition, we visualized the reconstruction error along the amount of errors injected.

Table IV shows the results for our reconstruction experiment. The second column shows the data set and the third column the amount of induced errors as percentage. As can be seen, the errors for the DS-AND are higher in comparison to the DS-SAN data set. This is due to the higher resolution

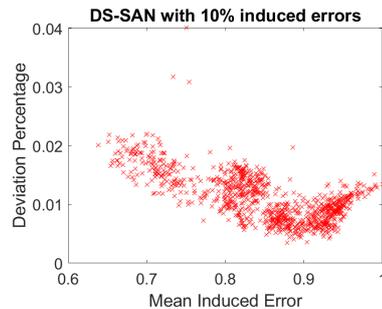


Fig. 5. The y axis represents the mean induced error scaled to a maximum of 1. The y axis represent the deviation percentage of the mean absolute error for the error correction in relation to the maximal mean induced error.

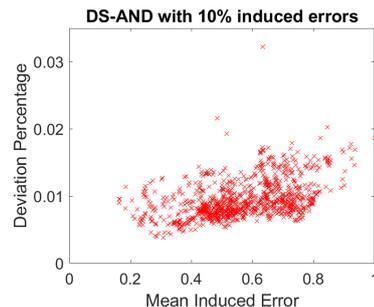


Fig. 6. The y axis represents the mean induced error scaled to a maximum of 1. The y axis represent the deviation percentage of the mean absolute error for the error correction in relation to the maximal mean induced error.

where the gaze points are mapped. The upper part shows the mean absolute error for reconstructing the entire input sequence. Since the neural network sees already a majority of values for reconstruction, the error is low. In contrast to this, the lower part of Table IV evaluates values only that where changed (Induced error). Of course the reconstruction error for those values increases, but it is interesting to see that for the data set DS-SAN the amount of induced error has only a slight impact in comparison to DS-AND. This is due to the sampling rate and that it is more likely to hit large position changes (Saccades) in DS-AND since Saccades in DS-SAN are only one or a few consecutive samples.

Figure 5 and Figure 6 show the mean absolute error for each induced error in a input sequence as a percentage to the maximum induced error (y axis). In addition, the x axis shows the induced error normalized to 1. Each error input sequence is represented as a red cross. As can be seen for both data sets is that the upper bound is at 2% with some outliers at 3% (DS-AND) and 4% (DS-SAN). The error for DS-AND behaves as you would expect. Meaning, that a larger error causes a larger reconstruction error. However, it is not the case for DS-SAN. Figure 5 clearly shows that the minimum reconstruction error is around 90% of the induced error. The reason for this is that the x axis represents the mean induced error of an input sequence and the y axis the mean reconstruction error for those errors. Since in DS-SAN it is less likely to hit a saccade, the mean induced error is most likely to be at 90%. This can also

be seen by the population of red crosses in Figure 5, which are most likely around these 90%. Therefore, our model learned to compensate for this linearly with the bias term of the neurons.

C. Evaluation Generation

Since we cannot directly compare generated scanpaths with other scanpaths, we decided to use a classification experiment. The data we used is from the ETRA 2019 challenge [43], [44]. It consists of 960 trials with a recording length of 45 seconds each. The recorded task are visual fixation, visual search and visual exploration. Since the visual fixation does not hold much complexity for generation, we omitted the data from those experiments. In addition four different stimuli were used in the experiments, which are: Blank, natural, where is waldo, and picture puzzle. For the classification itself, we used the approach proposed in [45]. This approach consists of transforming the eye tracking data into images. These images contain the raw gaze data as dots in the red channel, the time is encoded into the blue channel and the green channel contains the path as lines between raw gaze points. As classifier, we used the same network as proposed in [45].

The classification experiment consists of two parts. The first part uses our VAE to generate one thousand new examples for each stimulus since the exploration and search task are both marked as free viewing in the data set. Afterwards, we used the classification network to predict the stimuli on the generated data. For training, of the classifier we used 50% of the data set. The other 50% were used for training the generators, where each stimuli type was trained separately. For the generated scanpath, it has to be mentioned that they were centered on the image based on their mean value. This was done to avoid blank images and images where the scanpath is only partially drawn. For the second experiment, we used the VAE to generate data to improve the classification result. Therefore, we trained the generator on the same 50% and used the other 50% only for validation. This means that the generator and the classifier used the same real data for training where of course the classifier was also trained on additional 1,000 generated scanpath per stimuli. The generated scanpath was centered as mentioned before. For training of the classifier, we used the same parameters as in [45]. In addition, we used more advanced augmentation techniques. First, we added random noise by shifting a gaze point with a chance of 20% around 10% of its original location. The second augmentation technique was shifting the scanpath around 30% of its original central position (mean). Additionally, we used cropping of the input data which means that we used between 50-100% of a scanpath. Therefore, we selected the cropping length and starting index randomly. Table V shows the results for both experiments. The upper part (Real Data) is the evaluation of the classification on the test set. As can be seen, we achieved similar results as in [45]. For the first experiment, we want to evaluate or generated examples based on the classification. This is shown in the central part in Table V (Gen. Data). As can be seen, all stimuli achieved a classification accuracy above chance level (25%). This can

TABLE V
RESULTS FOR THE STIMULI CLASSIFICATION USING THE REAL DATA, THE GENERATED DATA AND THE GENERATED DATA ADDITIONALLY FOR TRAINING.

		Blank	Natural	Puzzle	Waldo	Accuracy
Real Data	Blank	17	12	0	3	0.531
	Natural	3	50	0	3	0.892
	Puzzle	1	1	58	0	0.966
	Waldo	0	6	2	52	0.866
Gen. Data	Blank	386	350	59	205	0.386
	Natural	267	422	149	162	0.422
	Puzzle	400	73	419	108	0.419
	Waldo	281	181	45	493	0.493
Gen. Train	Blank	21	11	0	0	0.656
	Natural	2	53	0	1	0.946
	Puzzle	0	0	60	0	1.0
	Waldo	0	2	1	57	0.95

be interpreted as our generated examples contain information about the gaze behavior from the specific stimuli. In addition, for each Stimuli, the second most classified target is Blank (For the true target Blank it is Natural). Those two observations mean that, either our generated data can be mapped to random gaze behavior (Blank means the screen only contains the gray color), or that it contains useful information that could not be learned from the training data so far. Therefore, we conducted our second experiment where we used additionally 4,000 generated examples for training. The results can be seen in the lower part in Table V (Gen. Train). As you can see by the results, the generated data is helpful in improving the classification results. One reason for this is that the model has to learn different combinations of gaze behavior and thus rather learns important patterns. This helps the model to generalize. In addition the data set is more balanced with the additionally generated data (Blank was underrepresented in the original data set).

V. CONCLUSION

We showed that, based on the input tensor construction, it is possible to use raw eye tracking data with fully convolutional neural networks for multiple tasks. They have the additional advantage that they can be used with any input size. In our results, we are improving the state of the art in the field of eye movement classification. Our main contribution in this area, however, is the construction of the input tensor as well as the pre-initialization of the first layer. This allows the use of raw data and makes this approach easy to use. In addition, the same approach can be used to improve data quality for experiments already performed. This is also a useful application as seen by the authors. Another interesting contribution of this work is the use of VAE for data generation. Compared to GANs, they are easier to train and can be combined with them for further improvement. Generating gaze data is also useful for testing many applications where the main purpose of course remains in the realm of training data generation.

REFERENCES

- [1] T. C. Kübler, C. Rothe, U. Schiefer, W. Rosenstiel, and E. Kasneci, "Subsmatch 2.0: Scanpath comparison and classification based on subsequence frequencies," *Behavior Research Methods*, vol. 49, no. 3, pp. 1048–1064, 2017.
- [2] J. G. May, R. S. Kennedy, M. C. Williams, W. P. Dunlap, and J. R. Branman, "Eye movement indices of mental workload," *Acta psychologica*, vol. 75, no. 1, pp. 75–89, 1990.
- [3] Q. Ji, Z. Zhu, and P. Lan, "Real-time nonintrusive monitoring and prediction of driver fatigue," *IEEE transactions on vehicular technology*, vol. 53, no. 4, pp. 1052–1068, 2004.
- [4] L. L. Di Stasi, M. B. McCamy, S. L. Macknik, J. A. Mankin, N. Hoof, A. Catena, and S. Martinez-Conde, "Saccadic eye movement metrics reflect surgical residents' fatigue," *Annals of surgery*, vol. 259, no. 4, pp. 824–829, 2014.
- [5] R. J. Leigh and D. S. Zee, *The neurology of eye movements*. Oxford University Press, USA, 2015, vol. 90.
- [6] A. Siekawa, M. Chwesiuk, R. Mantiuk, and R. Piórkowski, "Foveated ray tracing for vr headsets," in *International Conference on Multimedia Modeling*. Springer, 2019, pp. 106–117.
- [7] K. Holmqvist, M. Nyström, R. Andersson, R. Dewhurst, H. Jarodzka, and J. Van de Weijer, *Eye tracking: A comprehensive guide to methods and measures*. OUP Oxford, 2011.
- [8] R. Andersson, L. Larsson, K. Holmqvist, M. Stridh, and M. Nyström, "One algorithm to rule them all? an evaluation and discussion of ten eye movement event-detection algorithms," *Behavior Research Methods*, vol. 49, no. 2, pp. 616–637, 2017.
- [9] A. T. Duchowski and S. Jörg, "Modeling physiologically plausible eye rotations," *Proceedings of Computer Graphics International*, 2015.
- [10] W. Fuhl and E. Kasneci, "Eye movement velocity and gaze data generator for evaluation, robustness testing and assess of eye tracking software and visualization tools," *CoRR*, vol. abs/1808.09296, 2018. [Online]. Available: <http://arxiv.org/abs/1808.09296>
- [11] D. D. Salvucci and J. H. Goldberg, "Identifying fixations and saccades in eye-tracking protocols," in *Eye Tracking Research and Applications*. ACM, 2000, pp. 71–78.
- [12] H. Widdel, "Operational problems in analysing eye movements," *Advances in psychology*, vol. 22, pp. 21–29, 1984.
- [13] R. Engbert and R. Kliegl, "Microsaccades uncover the orientation of covert attention," *Vision research*, vol. 43, no. 9, pp. 1035–1045, 2003.
- [14] O. V. Komogortsev and J. I. Khan, "Eye movement prediction by oculomotor plant kalman filter with brainstem control," *Journal of Control Theory and Applications*, vol. 7, no. 1, pp. 14–22, 2009.
- [15] O. V. Komogortsev, D. V. Gobert, S. Jayarathna, D. H. Koh, and S. M. Gowda, "Standardization of automated analyses of oculomotor fixation and saccadic behaviors," *IEEE Transactions on Biomedical Engineering*, vol. 57, no. 11, pp. 2635–2645, 2010.
- [16] G. Veneri, P. Piu, P. Federighi, F. Rosini, A. Federico, and A. Rufa, "Eye fixations identification based on statistical analysis-case study," in *Cognitive Information Processing (CIP), 2010 2nd International Workshop on*. IEEE, 2010, pp. 446–451.
- [17] G. Veneri, P. Piu, F. Rosini, P. Federighi, A. Federico, and A. Rufa, "Automatic eye fixations identification based on analysis of variance and covariance," *Pattern Recognition Letters*, vol. 32, no. 13, pp. 1588–1593, 2011.
- [18] T. Santini, W. Fuhl, T. Kübler, and E. Kasneci, "Bayesian identification of fixations, saccades, and smooth pursuits," in *Eye Tracking Research and Applications*. ACM, 2016, pp. 163–170.
- [19] M. Nyström and K. Holmqvist, "An adaptive algorithm for fixation, saccade, and glissade detection in eyetracking data," *Behavior research methods*, vol. 42, no. 1, pp. 188–204, 2010.
- [20] R. van der Lans, M. Wedel, and R. Pieters, "Defining eye-fixation sequences across individuals and tasks: the binocular-individual threshold (bit) algorithm," *Behavior Research Methods*, vol. 43, no. 1, pp. 239–257, 2011.
- [21] L. Larsson, M. Nyström, and M. Stridh, "Detection of saccades and postsaccadic oscillations in the presence of smooth pursuit," *IEEE Transactions on Biomedical Engineering*, vol. 60, no. 9, pp. 2484–2493, 2013.
- [22] L. Larsson, M. Nyström, R. Andersson, and M. Stridh, "Detection of fixations and smooth pursuit movements in high-speed eye-tracking data," *Biomedical Signal Processing and Control*, vol. 18, pp. 145–152, 2015.
- [23] S. Hoppe and A. Bulling, "End-to-end eye movement detection using convolutional neural networks," *arXiv preprint arXiv:1609.02452*, 2016.
- [24] R. Zemblyns, D. C. Niehorster, O. Komogortsev, and K. Holmqvist, "Using machine learning to detect events in eye-tracking data," *Behavior research methods*, vol. 50, no. 1, pp. 160–181, 2018.
- [25] W. Fuhl, N. Castner, and E. Kasneci, "Rule-based learning for eye movement type detection," in *Proceedings of the Workshop on Modeling Cognitive Processes from Multimodal Data*. ACM, 2018, p. 9.
- [26] —, "Histogram of oriented velocities for eye movement detection," in *Proceedings of the Workshop on Modeling Cognitive Processes from Multimodal Data*. ACM, 2018, p. 5.
- [27] S. P. Lee, J. B. Badler, and N. I. Badler, "Eyes alive," in *ACM Transactions on Graphics (TOG)*, vol. 21, no. 3. ACM, 2002, pp. 637–644.
- [28] X. Ma and Z. Deng, "Natural eye motion synthesis by modeling gaze-head coupling," in *IEEE Virtual Reality Conference*. IEEE, 2009, pp. 143–150.
- [29] H. Murphy and A. T. Duchowski, "Perceptual gaze extent & level of detail in vr: looking outside the box," in *ACM SIGGRAPH conference Abstracts and Applications*. ACM, 2002, pp. 228–228.
- [30] B. H. Le, X. Ma, and Z. Deng, "Live speech driven head-and-eye motion generators," *IEEE transactions on Visualization and Computer Graphics*, vol. 18, no. 11, pp. 1902–1914, 2012.
- [31] D. Tweed, W. Cadera, and T. Vilis, "Computing three-dimensional eye position quaternions and eye velocity from search coil signals," *Vision research*, vol. 30, no. 1, pp. 97–110, 1990.
- [32] E. Wood, T. Baltrusaitis, X. Zhang, Y. Sugano, P. Robinson, and A. Bulling, "Rendering of eyes for eye-shape registration and gaze estimation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 3756–3764.
- [33] T. Pejsa, B. Mutlu, and M. Gleicher, "Stylized and performative gaze for character animation," in *Computer Graphics Forum*, vol. 32, no. 2pt2. Wiley Online Library, 2013, pp. 143–152.
- [34] D. J. Campbell, J. Chang, K. Chawarska, and F. Shic, "Saliency-based bayesian modeling of dynamic viewing of static scenes," in *Eye Tracking Research and Applications*. ACM, 2014, pp. 51–58.
- [35] A. Duchowski, S. Jörg, A. Lawson, T. Bolte, L. Świrski, and K. Krejtz, "Eye movement synthesis with 1/f pink noise," in *Proceedings of the 8th ACM SIGGRAPH Conference on Motion in Games*. ACM, 2015, pp. 47–56.
- [36] A. T. Duchowski, S. Jörg, T. N. Allen, I. Giannopoulos, and K. Krejtz, "Eye movement synthesis," in *Eye Tracking Research and Applications*. ACM, 2016, pp. 147–154.
- [37] D. Simon, S. Sridharan, S. Sah, R. Ptucha, C. Kanan, and R. Bailey, "Automatic scanpath generation with deep recurrent neural networks," in *Proceedings of the ACM Symposium on Applied Perception*. ACM, 2016, pp. 130–130.
- [38] M. Assens, X. Giro-i Nieto, K. McGuinness, and N. E. O'Connor, "Pathgan: visual scanpath prediction with generative adversarial networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 0–0.
- [39] H. Robbins and S. Monro, "A stochastic approximation method," *The annals of mathematical statistics*, pp. 400–407, 1951.
- [40] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014.
- [41] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [42] R. S. Hessels, D. C. Niehorster, C. Kemner, and I. T. Hooge, "Noise-robust fixation detection in eye movement data: Identification by two-means clustering (i2mc)," *Behavior research methods*, vol. 49, no. 5, pp. 1802–1823, 2017.
- [43] J. Otero-Millan, X. G. Troncoso, S. L. Macknik, I. Serrano-Pedraza, and S. Martinez-Conde, "Saccades and microsaccades during visual fixation, exploration, and search: foundations for a common saccadic generator," *Journal of vision*, vol. 8, no. 14, pp. 21–21, 2008.
- [44] M. B. McCamy, J. Otero-Millan, L. L. Di Stasi, S. L. Macknik, and S. Martinez-Conde, "Highly informative natural scene regions increase microsaccade production during visual scanning," *Journal of neuroscience*, vol. 34, no. 8, pp. 2956–2966, 2014.
- [45] W. Fuhl, E. Bozkir, B. Hosp, N. Castner, D. Geisler, T. C. Santini, and E. Kasneci, "Encodji: Encoding gaze data into emoji space for an amusing scanpath classification approach :)," in *Eye Tracking Research and Applications*, 2019.